

# Remote sensing image classification based on formal concept analysis

MAO Dian-hui<sup>1</sup>, LI Wen-zheng<sup>1</sup>, LIN Wei-hua<sup>2</sup>

1. School of Computer and Information Engineering, Beijing Technology and Business University, Beijing 100048, China;

2. School of Information Engineering, China University of Geoscience, Hubei Wuhan 430074, China

**Abstract:** In order to solve the problems that how to mine and express classification knowledge and rules in current remote sensing image classification, this paper introduces a new data mining theory of formal concept analysis, and realizes the connotation reduction of concept based on the minimum coverage of sets for ensuring the simplicity of classification rules. Meanwhile, the Fang city of Hubei province is selected to carry out the formal concept analysis theory to mine the land-use types classification rules, and construct a heuristic classifier based on the mined classification rules. The result shows that the mined classification rules have higher credibility, and the constructed classifier has higher accuracy compared with supervision classification and C4.5 algorithm, which proves that the theory of formal concept analysis provides a new method to achieve remote sensing image classification.

**Key words:** formal concept analysis, concept lattice, remote sensing image classification

**CLC number:** TP751.1      **Document code:** A

**Citation format:** Mao D H, Li W Z and Lin W H. 2010. Remote sensing image classification based on formal concept analysis. *Journal of Remote Sensing*. 14(1): 090—103

## 1 INTRODUCTION

Remote sensing image classification has been attached great importance by researchers over the years, which aims to confirm the discriminative criterions of different surface features (Sun, 1997). However, the disadvantages of traditional visual interpretation classification method were the poor timeliness and repeatability, and the interpretation results were difficult to compare and converse varying from person to person. In recent years, the automatic classification technology has been extensively carried out with the development of artificial intelligence, which included supervised classification, non-supervised classification, decision tree, neural network classification, support vector machine, expert knowledge classification and so on (Sun 2007). For example, Chen (1996) combined pattern recognition and artificial intelligence technology to classify the Landsat TM image; Xiong (2000) carried out the neural network algorithm on high spectral remote sensing imagery of Beijing; Chen (2007) utilized the decision tree algorithm C4.5 to mine land-use types classification rules on the Landsat ETM image of Fuzhou; He (2006) implemented support vector machine on the remote sensing image classification. The methods mentioned above improved the classification precision in some extent. However, the disadvantages of these methods in practi-

cal application included: the classification rules mined by neural network algorithm were difficult to understand because the knowledge were implicatively embedded in its network; the kernel function of support vector machine was hard to selected in experiment, and it has problem of classification knowledge expression which is the same the neural network algorithm; the algorithm of decision tree ignored the relations of attribute sets. Thus, in order to solve the deficiencies mentioned above, this paper introduces a new data mine theory of formal concept analysis (FCA) into remote sensing image classification, the theory can completely express the various modes among data attribute sets and provide a new idea to achieve the classification knowledge acquisition in remote sensing image classification research.

## 2 FORMAL CONCEPT ANALYSIS THEORY

The formal concept analysis theory proposed by Wille in 1982, was a theory of data analysis and rule mined according to building concept lattices based on the formal context (Wille, 1992). Every concept comprises intension and extension in concept lattice. The extension of concept represents a set of objects. The intension of concept represents the common features that all objects in the extension have. Concept lattice reflects entity-attribute relationships between objects. The corre-

**Received:** 2008-09-25; **Accepted:** 2008-12-08

**Foundation:** National Natural Science Foundation (No. 40801161/D010701).

**First author biography:** MAO Dian-hui (1979— ), male, PH.D, Lecture, Graduated from Huazhong University of science and technology. Majoring in pattern recognition and spatial data mining. He has published 11 papers. E-mail: amaode@gmail.com

(C)1994-2021 China Academic Journal Electronic Publishing House. All rights reserved. <http://www.cnki.net>

sponding Hasse diagram vividly visualizes the relations between concepts (Wille, 1989). Thus FCA has been adopted as a tools to analysis data and mine various knowledge from database.

2.1 Definition

**Definition 1:** A formal content is briefly defined as a triplet set  $K=(G,M,R)$ , where  $G$  is a set of objects,  $M$  is a set of attributes, and  $R$  is the binary relationships between  $G$  and  $M$ .

**Definition 2:** In the formal content  $K$ , there is a mapping relation between the power sets of  $G$  and  $M$ :

$$\forall A \subseteq G : f(A) = \{m \in M \mid \forall g \in A, gRm\}$$
$$\forall B \subseteq M : g(B) = \{g \in G \mid \forall m \in B, gRm\}$$

where the mapping relations of function  $f$  and  $g$  is called Galois mapping between the power sets of  $G$  and  $M$ , so it is also called as Galois lattice.

**Definition 3:** A couple  $(A, B)$  derived from formal context  $K$  is called a basic concept, which satisfies:

$$A \subseteq G, B \subseteq M, f(A) = B, g(B) = A$$

where  $A$  is called the extension of concept  $(A, B)$ , and  $B$  is called the intension of concept  $(A, B)$ .

**Definition 4:** In the concept lattice, partial order relation “ $<$ ” between concept  $C_1=(A_1, B_1)$  and  $C_2=(A_2, B_2)$  is defined as  $C_1 < C_2$ , when satisfies  $A_1 \subseteq A_2$  ( $\Leftrightarrow B_1 \supseteq B_2$ ).  $C_1$  is called the sub-concept (son) of  $C_2$ , and  $C_2$  is called the sup-concept (father) of  $C_1$ . That is, the relationship between  $C_2$  and  $C_1$  is the relationship of father and son. If  $C_1 < C_2$ , there is not a concept  $C=(A, B)$  that satisfies  $C_1 < C < C_2$ .  $C_1 < C_2$  is called an immediate-sub-concept-relation between  $C_1$  and  $C_2$ .

According to the partial order relation, concept lattice can be denoted by a labeled line diagram. Every node of the diagram represents a concept, and the line connecting the nodes expresses the relationship of generalization and specialization between these nodes. The line diagram is called Hasse diagram which is a visual denotation of concept lattice. Table 1 is the formal context  $K=(G, M, I)$ , and the Hasse diagram of its concept lattice is denoted by Fig. 1.

2.2 Construction algorithm of Concept Lattice

The process of constructing concept lattice is clustering concepts. There is only structure to the same data set by carrying out different algorithms to build concept lattice. Therefore, the merit of concept lattice is not affected by order of data or

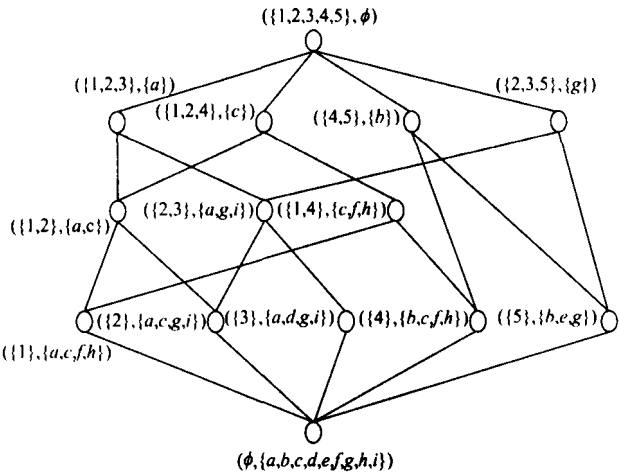


Fig. 1 Concept lattice of the example

attribute sets. On the other hand, the size of concept lattice is exponential times growth with the number of objects and properties of formal content. Thus, how to construct concept lattice is always a major research field. In this study, an efficient algorithm is carried out by dividing into two main processes: (1) to a known concept  $(A, B)$ , all direct super-concepts of which are calculated by Neighbors algorithm (Jian & Hu, 2001); (2) iteratively employing Neighbors algorithm to every concept to generate Hasse diagram. The algorithm can quickly estimate whether a node of Hasse diagram has been generated or not according to organizing the generated nodes as a dictionary index.

2.3 Method of mining classification rules

Known from the characteristics of concept lattices, the intension of every concept is viewed as a maximal frequent item-set. Therefore, concept lattice as one of data structures is very suitable to extract association rules from database, because the classification rule is a special form of association rule when the latter parts of rule are the decision-making properties of rules. Thus, to a concept who included sort property, we only need to judge whether the node satisfied the threshold of confidence and support or not. If the concept met the conditions, the classification rule is expressed that the sort property is regarded as the latter part of classification rules, and other attributes of the concept are viewed as the condition parts of the rule.

**Theorem:** Supposing  $K = \{G, M, R\}$  is a decision-making content,  $M = C \cup D$ ,  $R = R_C \cup R_D$ , where  $C$  is the condition attribute set,  $D$  is the sort attribute set,  $R_C \subseteq G \times C$  is the condition relation set, and  $R_D \subseteq G \times D$  is the attribute relation set. If  $B_{C1} \subseteq C, B_{D1} \subseteq D$ , then the rule  $B_{C1} \rightarrow B_{D1}$  of  $K$  is a classification rule when it satisfies  $B_{C1}^C \subseteq B_{D1}^D$  (Zhai, 2006).

The support degree of classification rule is defined as  $S(B_{C1} \rightarrow B_{D1}) = \frac{N_{ins\_class}(B_{C1}, B_{D1})}{|U|}$ , where  $N_{ins\_class}(B_{C1}, B_{D1})$  is the number of instances which belong to class  $B_{D1}$  and

Table 1 Example of formal context

		M								
G	I	a	b	c	d	e	f	g	h	i
	1	1	0	1	0	0	1	0	1	0
	2	1	0	1	0	0	0	1	0	1
	3	1	0	0	1	0	0	1	0	1
	4	0	1	1	0	0	1	0	1	0
	5	0	1	0	0	0	0	1	0	0

meet the condition attribute set  $B_{C1}$  (consider the repetitious instances), and  $U$  is the number of all instances in the formal content.

The calculated method of support degree are as follows: firstly, counting the number of the extension set of the top node in the concept lattice, which includes all objects in formal content; secondly, counting the number of the extension set of the nodes which include sort attributes. Therefore, the formula of support degree of the rule is:  $\text{Sup} = \frac{O_D}{O_{\text{top}}} \times 100\%$ , where  $O_D$  is

the number of the extension set of the nodes who fit with the qualification, and  $O_{\text{top}}$  is the extension set of the top concept.

The confidence degree of the rule is defined as:

$$c(B_{C1} \rightarrow B_{D1}) = \frac{N_{\text{ins\_class}}(B_{C1}, B_{D1})}{N_{\text{ins\_ref}}(B_{C1})}, \text{ where } N_{\text{ins\_ref}}(B_{C1}) \text{ is the}$$

number of instances who satisfy the condition attribute set  $B_{C1}$  (consider the repetitious instances),  $N_{\text{ins\_ref}}(B_{C1}, B_{D1})$  is the number of instances who belong to class  $B_{D1}$  and meet the condition attribute set  $B_{C1}$  (consider the repetitious instances).

The calculated method of the confidence degree is as follows: to a concept who includes sort attribute, firstly, counting the number of the extension set of its parent node and judging whether the intension set of the parent node is composed of the condition attribute set of the original concept node or not. If the parent concept node exists, then counting the number of the extension set of the nodes. Therefore, the formula of confidence

degree is:  $\text{Conf} = \frac{O_D}{O_{\text{parent}}} \times 100\%$ , where  $O_D$  is the number of

the extension set of the nodes which fit with the qualification, and  $O_{\text{parent}}$  is the number of the extension set of the parent concept. If the parent concept node does not exist, then the confidence degree is 100%.

## 2.4 Intension reduction of concept

In order to guarantee the property set of  $B_{C1} \rightarrow B_{D1}$  without redundancy in the formal content, the concept nodes required to reduce intension attributes while maintaining the extension set without change, which is defined as:

**Definition 5:** To a given concept  $C = (O_1, D_1)$ , if the set  $D_2$  satisfies the conditions:

- (1)  $g(D_2) = g(D_1) = O_1$ ;
- (2) To arbitrary  $D_3 \subset D_2$ , exists  $g(D_3) \supset g(D_2) = O_1$ ;

Then the set  $D_2$  is viewed as one of the reduction sets of intension of the concept  $C$ .

In order to fulfill with the intension reduction of concept, the minimum coverage theory of group sets is carried out in this paper as follows:

**Definition 6:** To a given group sets  $FM = \{M_1, M_2, \dots, M_n\}$ , where the set  $M$  is defined as the minimum coverage of  $FM$ , if it meets the condition:

- (1)  $\forall M_i \in FM (M \cap M_i \neq \emptyset)$ ;
- (2)  $\forall M' \subset M (\exists M_i \in FM (M' \cap M_i \neq \emptyset))$ ;

**Definition 7:** To a given concept node  $C = (O_1, D_1)$  and its subset  $D_2 \subseteq D_1$ ,  $D_2$  is the minimum coverage of group sets  $\{D_1 - D_3 | (O_3, D_3) \text{ is the parent concept of } C\}$ , where  $D_2$  is defined as one of the reduction sets of intension of the concept  $C$ .

According to the definition mentioned above, how to dispose the intension reduction of concept is translated into how to get the minimum coverage of group sets, the algorithm of how to get the minimum coverage of group sets is described in detail in reference (Xie, 2001).

## 3 EXPERIMENT AND ANALYSIS

In order to carry out the formal concept analysis theory on remote sensing image classification, the Fang City of Hubei province is selected as the study area in this paper. The collected original data include Landsat ETM (acquired on June 15, 2000), 1 : 100000 digital topographic maps, administrative divisions maps, land-use types maps in 2000.

### 3.1 Feature extraction of remote sensing images

The spectral characteristics value of different bands is the corresponding DN values of bands (band TM6 is not chosen because of its lower-resolution compared with others in this paper), and the correlation characteristics of bands are denoted as the normalized difference index between bands. The formula

$$\text{is: } \text{NDI}_{ij} = \frac{\text{TM}_i - \text{TM}_j}{\text{TM}_i + \text{TM}_j} (i, j = 1, 2, 3, 4, 5, 7; i \neq j). \text{ The texture}$$

characteristics are extracted by carrying out grey-level co-occurrence matrix method on the panchromatic band of Landsat ETM, in which the moving window of the method is set as 3×3, the parameter of moving length is set as 1, and the value of moving angle is set as 45°. The selected eight texture characteristics include mean, variance, inverse variance, contrast, non-similarity, entropy, second moment and relevance.

To the geographic data, 1 : 100000 topographic maps are firstly vectorized to DEM. Then slope maps are extracted from DEM on the Arcgis 9.2 platform according to re-sampling and spatial-adjust referring to the original remote sensing images. Lastly, the elevation and slope features data are projected as two new "bands" to the uniform coordinate system of remote sensing images.

### 3.2 Data processing and analysis

In order to reduce errors caused by artificial selecting sample data, an auto image clipping program is developed by IDL based on the ENVI 4.2 platform. The unit of image clipping object is the plot of land-use types map. Meanwhile in order to avoid the data error caused by interpreting accuracy of land-use types map, the most frequency value of histogram in plots is statistic as the characteristic values, which compose the multi-resource spatial database. The database structure is shown as Fig. 2.

According to field survey and data statistics in study area, the main land-use types were divided into grassland, woodland, dry land, paddy fields and water body. Meanwhile, in order to guarantee the accuracy of classification rules mined from the database, 5133 sample data are selected and every land-use type approximately equals. On the other hand, convenient for comparative analysis of the various characteristic values, every characteristic value is normalized to range 0—255 by the for-

mulas:  $Y = \frac{y - y_{\min}}{y_{\max} - y_{\min}} \times 255$ , where the range of normalized

difference index between bands are  $(-1, 1)$ , the scope of elevation are  $(200\text{m}, 2200\text{m})$ , and the bound of slope are  $(0^\circ, 69^\circ)$ . The statistical curves of all characteristics are shown in Fig. 3.

Known from the curves of spectral character Fig. 3(a), the

plots of different land-use types are obviously distinguished in TM3, but not clear in TM1 and TM5. Meanwhile, water-body can be distinctly discerned from others except TM5, but the DN value of dry land, grassland, paddy field and woodland are relatively similar in all bands. However, it is not easy to distinguish dry land and paddy field from grassland and woodland because the selected images are obtained in summer and crops grow well.

Fig.3 (b) shows the curves of normalized difference index between bands (NDI). Abscissas 1—15 represent NDI12, NDI13, NDI14, NDI15, NDI17, NDI23, NDI24, NDI25, NDI27, NDI34, NDI35, NDI37, NDI45, NDI47, and NDI57, respectively. The curves show that the NDI of TM3 and TM4, TM7 and TM4, TM3 and TM5, TM1 and TM3, TM1 and TM7

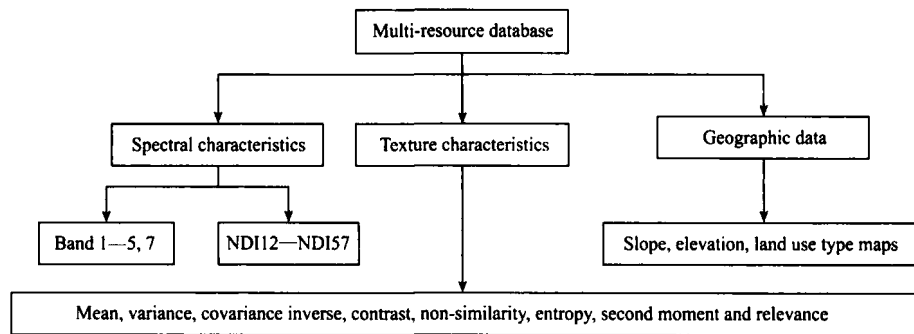


Fig. 2 Structure of multi-resource database

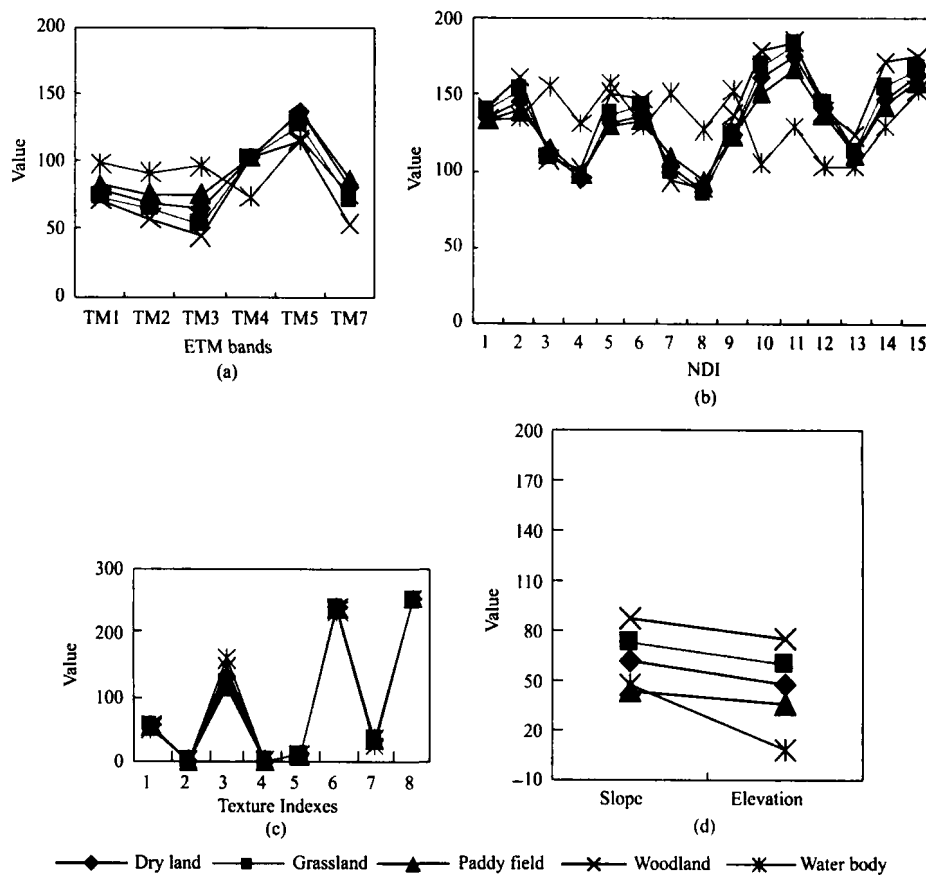


Fig. 3 Curves of different land use type characters

(a) Spectral characters; (b) Normalized difference between bands; (c) Texture characters; (d) Slope and elevation characters

have higher ability to distinguish, while the index between TM1 and TM2 was almost unable to discern. Therefore, the NDI of TM1 and TM2 are not selected for mining classification rules.

Fig. 3 (c) shows the curves of texture characters. Abscissas 1—8 represent mean, variance, inverse variance, contrast, non-similarity, entropy, second moment, relevance, respectively. Known from the curves of texture characters, all characters are almost entirely consistent except inverse variance, so the index of inverse variance is only chosen in this study.

Known from the curves of elevation and gradient, paddy fields and water-body are mainly distributed in plains, but dry land is mostly located in hilly area, and woodland or grassland are generally distributed in mountains due to less man-made factors.

3.3 Classification rules extraction and analysis

For concept lattice can only dispose discrete data, every

characteristic index is divided into several discrete intervals according to data distribution histogram in this study. So every discrete interval is expressed as: Name\_XX, where Name is denoted as the character index, XX is the number of the character interval. Spectral characters, normalized difference index between bands, as well as the texture characters are divided into 25 intervals, but gradient is divided into seven intervals in accordance with the national standard as shown in Table 2, and elevation is divided into 12 intervals for that the plots can be discerned where the distribution gap is about 120 m in vertical terrain.

For the discrete sample data, 9033 concepts are generated based on the constructed algorithm of concept lattice in this study, in which 4835 concept nodes contain sort attribute. And then the concepts disposed by the intension reduction algorithm generate corresponding classification rules, only partly one-dimensional classification rules of higher support degree are shown in Table 3 due to space limitations.

Table 2 Intervals of slope

0°—3°	3°—5°	5°—8°	8°—15°	15°—25°	25°—35°	>35°
Slope_1	Slope_2	Slope_3	Slope_4	Slope_5	Slope_6	Slope_7

Table 3 Single dimension attribute classification rules

No	Rules	Support /%	Confidence /%
1	band2_6—>dry land	9.71	41.88388
2	ND47_16—>grassland	7.36	44.59759
3	band6_6—> grassland	6.49	40.80674
4	ND23_15—> grassland	6.55	50.1248
5	ND17_11—> paddy fields	7.51	50.55545
6	ND47_11—> paddy fields	8.16	46.58692
7	ND47_14—> paddy fields	4.33	55.43247
8	band4_8 —> paddy fields	5.25	44.76058
9	ND35_14 —> paddy fields	4.59	55.01782
10	ND13_12 —> paddy fields	3.42	50.21281
11	ND34_11 —> paddy fields	3.25	47.71157
12	ND35_15—> paddy fields	5.16	49.90062
13	ND47_18—> woodland	8.97	42.78857
14	ND35_17—> woodland	7.35	44.75626
15	ND27_13—> woodland	5.15	55.83393
16	dem_4—> woodland	5.71	53.49792
17	slope_3—> water body	6.71	43.30056
18	ND14_21—> water body	7.45	44.1342
19	ND27_20—> water body	4.70	42.42167
20	band4_1 —> water body	5.44	48.22366
21	dem_0 —> water body	16.30	48.21919

Known from the mined classification rules, multi-dimensional classification rules are generally composed of several single dimensional rules. However, the single dimensional rules have higher support degree, that is to say, they have more ability to distinguish in essence.

The conclusion analyzed from the mined one dimensional classification rules showed that: the rule of No.16 can be understood that when the elevation is about 827—984m, the support degree of the plots labeled as the woodland is up to 5.71%, and the confidence degree is up to 53.50 %; the rule of No.21 can be interpreted as that when the elevation is about 200—307m, the support degree of plots classified as water body is up to 16.30%, and the confidence degree is up to 48.2%; the result accords with the analysis of elevation and gradient, where water body is located in relatively flat plains and hills, but forests are distributed in higher mountains, so the mined rules can be understood and accepted by man. To the rules of No.2,6,7 and 13, the value of ND47\_18 of woodland > the value of ND47\_16 of grassland > the value of ND47\_11-ND47\_14 of paddy fields, which are coincident with the analysis of normalized difference index between bands. Therefore the result shows that all mined classification rules have higher credibility. On the other hand, comparing the length of mined classification rules, the result shows that the length of rules mined by concept lattice is generally about 3—6, less than the length of rules obtained by decision tree algorithm C4.5 where the length is average 8—10 attributes (Chen, 2007). So the rules mined by formal concept analysis are relatively brief, and the classifier built based on the rules will also be relatively simple.

3.4 Classifier construction and analysis

In order to construct a classifier based on the mined classification rules, a heuristic classifier is carried out in this study, the process included: firstly, all conflicting rules in rules set are

deleted; secondly, the rules set order in partial relationship. The order principles is: (1) confidence degree priority; (2) when confidence degree equals, the support degree priority; (3) when confidence and support degree both equal, the rules sort by generating order; (4), to the ordered rules set, constructing a classifier by selecting the subset of rules which has least train errors on the sample data. The algorithm of constructing heuristic classifier is described in detail by Qi (2005). Meanwhile, in order to study the classifier accuracy in different support and confidence threshold as well as evaluate the classifier's performance, some training sample data are chosen from multi-resource spatial databases, and carrying out experiment in different support and confidence thresholds. The result shows that when the threshold of support degree is 2% and confidence threshold is 40%, the heuristic classification accuracy is the highest and reaches 80.26%, which is 78.64% higher than that obtained by C4.5 algorithm based on the same sample data set.

In order to verify and compare the accuracy of heuristic classifier constructed by the classification rules based on the theory of the concept lattice, the remote sensing image orbited on 126/38 of May 5, 2000 is selected for experiment in this paper. 1166 sample plots are chosen, which include 275 dry land, 107 grass land, 236 paddy fields, 228 water body and 320 woodland. The results of confusion matrix method compared with supervised classification and decision tree classification method are shown in Table 4—Table 6.

Conclusions can be made from the above results that the total classification accuracy and Kappa coefficient of heuristic classifier are relatively higher compared with C4.5 decision tree algorithm and supervised classification method. The user's accuracy and producer's accuracy also shows some advantages, especially in distinguishing dry land, grasslands, paddy fields, forest land compared with supervision classification method, but slightly decreased in distinguishing water body compared

Table 4 Result of supervised classification

Type	Dry land	Grass land	Paddy fields	Water body	Woodland	Total	User's accuracy/%
Dry land	204	32	25	0	14	275	74.18
Grass land	19	62	17	0	9	107	57.94
Paddy fields	33	14	169	8	12	236	71.61
Water body	0	0	25	203	0	228	89.04
Woodland	23	38	25	9	225	320	70.31
Total	279	146	261	220	260	1166	—
Producer's accuracy/%	73.12	42.47	64.75	92.27	86.54	—	—

Total accuracy: 74.01%    Kappa coefficient: 0.7389

Table 5 Results of C4.5 algorithm

Type	Dry land	Grass land	Paddy fields	Water body	Woodland	Total	User's accuracy/%
Dry land	215	21	28	1	10	275	78.18
Grass land	21	71	8	0	7	107	66.35
Paddy fields	35	9	173	4	15	236	73.31
Water body	2	5	8	205	8	228	89.91
Woodland	20	42	19	5	234	320	73.13
Total	293	148	236	215	274	1166	—
Producer's accuracy/%	73.38	47.97	73.31	95.35	85.40	—	—

Total accuracy: 77.02%    Kappa coefficient: 0.7691

Table 6 Results of constructed heuristic classifier

Type	Dry land	Grass land	Paddy fields	Water body	Woodland	Total	User's accuracy/%
Dry land	221	25	14	2	13	275	80.36
Grass land	23	75	3	1	5	107	70.09
Paddy fields	27	12	181	2	14	236	76.69
Water body	4	6	15	199	4	228	87.28
Woodland	18	37	17	6	242	320	75.63
Total	293	155	230	210	278	1166	—
Producer's accuracy/%	75.43	48.39	78.70	94.76	87.05	—	—

Total accuracy: 78.73%    Kappa coefficient: 0.7862

with the decision tree algorithm C4.5. Although observational data is absent and experimental data selected from interpreted land-use types exist error, which result that the classification accuracy of three methods are relatively lower, the results compared among the three methods show that heuristic classifier has obviously advantage, which proved that the theory of formal concept analysis provides a new way to solve the problem of remote sensing image classification.

#### 4 CONCLUSIONS AND OUTLOOK

According to the formal concept analysis theory in this paper, the classification rules are mined from remote sensing image and a heuristic classifier is constructed based on the mined rules. The experimental result shows that the method has demonstrated advantage compared with supervised classification and decision tree algorithm C4.5, which also proved that the method provides a new way to mine classification rules. However, the shortcoming of formal concept analysis theory is sensitive to noise data in mining rules. When dealing with mixed areas in remote sensing image, how to decompose the mix pixels by formal concept analysis theory need to be further researched.

#### REFERENCES

- Chen X Y. 2007. The Classification of Remote Sensing Image Based on Spatial Data Mining and Knowledge Discovery. Fujian: Fujian Normal University
- Chen Y F, Wang Z Q, Wang Y G and Zhang Y. 1996. Expert system applied for TM image classification. *Forest Research*, 4(9): 344—347
- He L M. 2006. Support Vector Machines Ensemble and Its Application in Remote Sensing Classification. Zhejiang: Zhejiang University
- Jian S Q, Hu X G and Jiang M H. 2001. Incremental algorithms of extended concept lattice. *Computer Engineering and Applications*, 15 (37): 132—134
- Qi H. 2005. Research on Knowledge Discovery based on Lattice Concept Analysis. Jiling: Jiling University
- Sun J B, Shu N and Guan Z Q. 1997. The Remote Sensing Theory, Method and Application. Beijing: Surveying and Mapping Press
- Sun X B, Fan W, Yan P, Huang Y and Ma Y H. 2007. A brief review on the classification methods of land cover based on remote sensing image. *Chinese Agricultural Science Bulletin*, 9(23): 607—610
- Wille R. 1989. Knowledge Acquisition by Methods of Formal Concept Analysis, Learning Symbolic and Numeric Knowledge. New York: Nova Science Publishers, Inc
- Wille R. 1982. Restructuring lattice theory: an approach based on hierarchies of concepts. *Ordered Sets*. Dordrecht: D Reidel Publishing Company
- Xie Z P. 2001. Research on Knowledge Discovery Based on Concept Lattice Model. Hefei: Hefei University of Technology
- Xiong Z, Dong Q X and Zhen L F. 2000. High-rank artificial neural network algorithm for classification of hyper spectral image data. *Journal of Image and Graphics*, 3(5): 196—201
- Zhai Y H. 2006. Knowledge Representation and Knowledge Acquisition on Formal Concept Analysis. Shanxi: Shanxi University

# 基于形式概念分析的遥感影像分类

毛典辉<sup>1</sup>, 李文正<sup>1</sup>, 林伟华<sup>2</sup>

1. 北京工商大学 计算机与信息工程学院 北京 100048;

2. 中国地质大学信息工程学院 湖北 武汉 430074

**摘要：**针对目前遥感影像分类方法中存在分类知识难以获取的不足，尝试引入形式概念分析的数据挖掘理论，并基于族集最小覆盖理论实现概念内涵的缩减，从而保证分类规则的简洁与无冗余性。研究选取湖北省房县作为试验区，实现了该理论在研究区中土地利用类型分类规则的挖掘应用。基于挖掘出的分类规则构建了启发式分类器，实验结果表明形式概念分析理论挖掘出的分类规则可信度较高，基于挖掘出的分类规则构建的分类器相对于监督分类方法、决策树 C4.5 算法在分类精度上有一定优势，从而证明了它对遥感影像分类提供一种的新方法。

**关键词：**形式概念分析，概念格，遥感影像分类

**中图分类号：**TP751.1

**文献标识码：**A

**引用格式：**毛典辉, 李文正, 林伟华. 2010. 基于形式概念分析的遥感影像分类. 遥感学报, 14(1): 090—103

Mao D H, Li W Z and Lin W H. 2010. Remote sensing image classification based on formal concept analysis. *Journal of Remote Sensing*. 14(1): 090—103

## 1 引言

遥感影像分类作为遥感技术的重要分支，多年来一直受到研究人员的普遍重视，其核心任务就是确定不同地物类别间的判别界面和判别准则(孙家柄等, 1997)。传统的分类方法大多采用目视解译法，其时效性、可重复性差，解译结果也因人而异，很难进行比较和转换。近年来，随着人工智能技术的发展，在遥感影像计算机自动分类领域中开展了大量的研究，其研究方法主要包括：监督与非监督分类方法、决策树分类法、神经网络分类、支持向量机分类以及专家知识分类方法等(孙学邦等, 2007)。如陈永富等(1996)采用模式识别与人工智能技术相结合的专家系统分类对 TM 遥感图像进行分类研究；熊桢、童庆禧等(2000)利用高阶神经网络算法对北京市沙河镇地区的高光谱数据进行了分类实验；陈小瑜(2007)利用决策树算法 C4.5 对福州市的 Landsat ETM 影像提取土地利用类型分类规则等；何灵敏(2006)等利用支持向量机理论实现对遥感图像的分

类研究。虽然这些方法在一定程度上提高了分类精度，取得了较好的效果，但是由于分类方法自身特点在实际应用中还存在不足：如神经网络分类方法其分类知识蕴含在网络中，分类规则难以理解；支持向量机分类中核函数的选取难题以及与神经网络分类方法类似的知识表达问题；决策树分类虽然可以显式表达分类规则，但是忽略了数据集中属性之间的相关性。为了解决以上方法的不足，本文尝试在遥感影像分类领域中引入一种能够完备表达数据中各种模式的新型数据挖掘理论——形式概念分析理论，为实现分类知识的获取与表达提供一种新的思路与方法。

## 2 形式概念分析理论

形式概念分析(Formal Concept Analysis, 简称 FCA)是德国 Wille 教授(1982)根据哲学中概念思想提出的一种从形式背景建立概念格进行数据分析和规则提取的理论。其主要思想为概念格中的每个节点表示一个形式概念，每个概念由 2 部分组成：外

收稿日期：2008-09-25；修订日期：2008-12-08

基金项目：国家自然科学基金(编号：40801161/D010701)。

第一作者简介：毛典辉(1979—)，男，博士，讲师，2008年毕业于华中科技大学系统分析与集成专业，主要从事模式识别、空间信息处理等研究工作，发表论文 11 篇，E-mail: amaode@gmail.com。



延和内涵。其中概念的外延代表一组对象，内涵则为这组对象所具有的公共特征(属性)(Wille,1989)，与概念格相对应的 Hasse 图则形象地揭示了概念间的泛化和例化关系，实现了对数据的可视化，非常适用于从数据库中进行知识发现的描述，从而成为数据分析和规则提取的一种有效工具。

2.1 基本定义

定义 1: 一个形式背景定义为一个三元组  $K=(G, M, R)$ ，其中  $G$  是对象集合， $M$  是属性集合， $R \subseteq G \times M$  是  $G$  与  $M$  之间的一个二元关系。

定义 2: 在形式背景  $K$  中，在  $G$  的幂集和  $M$  的幂集之间可以定义两个映射  $f$  和  $g$  如下：

$$\forall A \subseteq G : f(A) = \{m \in M \mid \forall g \in A, gRm\}$$

$$\forall B \subseteq M : g(B) = \{g \in G \mid \forall m \in B, gRm\}$$

式中， $f$  和  $g$  被称为  $G$  的幂集和  $M$  的幂集之间的 Galois 连接。

定义 3: 形式背景  $K=(G, M, R)$  上的一个形式概念定义为一个二元组  $(A, B)$ ，满足：

$$A \subseteq G, B \subseteq M, f(A) = B, g(B) = A$$

式中， $A$  称为概念  $(A, B)$  的外延， $B$  称为概念  $(A, B)$  的内涵。

每一个概念描述了一组对象及其公共的特征。概念的内涵是概念外延中所有对象的共同属性的集合，概念的外延是概念内涵可以确定的最大对象的集合。

定义 4: 概念  $C_1=(A_1, B_1)$  和  $C_2=(A_2, B_2)$ ，如果  $A_1 \subseteq A_2 (\Leftrightarrow B_1 \supseteq B_2)$ ，则称概念  $C_1$  是概念  $C_2$  的子概念(也称后继)，概念  $C_2$  是概念  $C_1$  的超概念(也称为前驱)，记为  $C_1 < C_2$ 。

通过这种序关系，得到一个有序集  $(B(K), \leq)$ ，称为形式背景  $K$  的概念格。概念格是所有形式概念在子概念和超概念下的序集。因此概念格可通过 Hasse 图图形化表示，这使得给定数据背景的概念结构变得清晰和易于理解，从而实现概念格的可视化显示。如表 1 表示一个形式化背景实例，图 1 为其对应的概念格结构。

表 1 形式背景例子

		M								
G	I	a	b	c	d	e	f	g	h	i
	1	1	0	1	0	0	1	0	1	0
	2	1	0	1	0	0	0	1	0	1
	3	1	0	0	1	0	0	1	0	1
	4	0	1	1	0	0	1	0	1	0
	5	0	1	0	0	0	0	1	0	0

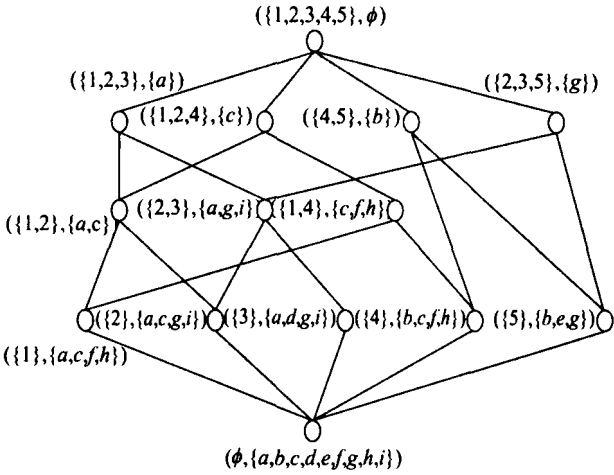


图 1 形式背景例子对应的概念格

2.2 概念格的构造

概念格的构造过程实际上是概念聚类过程，对于同一形式背景，采用不同的建格算法所构造的概念格是唯一的，因此概念格具有不受数据或属性排列次序影响的优点，但是另一方面，概念格的规模大小会以形式背景中对象个数和属性个数的指数倍增长，因此概念格构造算法是该领域研究的一个主要问题。在本研究中采用了高效的 Lattice 算法(简宋全等, 2001)，该算法主要分为 2 个过程：(1)已知概念  $(A, B)$ ，使用算法 Neighbors 计算出概念节点  $(A, B)$  的所有直接超概念节点集合；(2)对于每个概念  $(A_i, B_i)$ ，算法迭代调用 Neighbors，去生成概念格的 Hasse 图结构，此算法将所有已经生成的概念格节点通过一棵词典排序树组织起来作为一个索引结构，从而可以快速的判断某个节点是否已经生成。

2.3 分类规则提取方法

从概念格的特点可知，概念格中每个节点的内涵实际上就是一个最大项目频繁集，因此，概念格非常适合作为关联规则提取的数据结构。分类规则作为关联规则的一种特殊形式，当其规则后件为决策属性的某个类别时，得到的规则便为分类规则。因此，基于概念格结构求取分类规则，只需提取格结构中包含类属性的概念节点，并计算该节点是否满足用户定义的置信度与支持度阈值，对于满足条件的概念节点，将概念内涵中的类属性作为规则的后件决策属性，其余属性作为规则的前件条件属性。

定理：设  $K=\{G, M, R\}$  为决策背景， $M = C \cup D$ ， $R = R_c \cup R_d$ ，其中  $C$  为条件属性， $D$  为决策属性， $R_c \subseteq G \times C$  为条件关系集合， $R_d \subseteq G \times D$  为决策关

系集合。若  $B_{C1} \subseteq C$ ,  $B_{D1} \subseteq D$ , 则  $K$  上的规则  $B_{C1} \rightarrow B_{D1}$  为决策规则且仅当  $B_{C1}^C \subseteq B_{D1}^D$  (翟岩慧, 2006)。

规则的支持度定义为:  $S(B_{C1} \rightarrow B_{D1}) = \frac{N_{ins\_class}(B_{C1}, B_{D1})}{|U|}$ , 其中  $N_{ins\_class}(B_{C1}, B_{D1})$  是指满足条件属性  $B_{C1}$  且属于类  $B_{D1}$  的实例数目(考虑重复实例),  $U$  为形式背景中所有实例。其计算方法为: 首先取得概念格中顶层概念节点的外延基数, 因为该节点外延对应于所有对象; 其次对于每一个包含类属性的概念节点, 求取该节点的外延基数, 因此, 所得规则的支持度为:  $Sup = \frac{O_D}{O_{top}} \times 100\%$ , 其中  $O_D$  为待求节点的外延基数,  $O_{top}$  为顶层节点的外延基数。

规则的置信度定义为  $c(B_{C1} \rightarrow B_{D1}) = \frac{N_{ins\_class}(B_{C1}, B_{D1})}{N_{ins\_ref}(B_{C1})}$ , 其中  $N_{ins\_ref}(B_{C1})$  是指满足条件属性  $B_{C1}$  的实例数目(考虑重复实例),  $N_{ins\_class}(B_{C1}, B_{D1})$  为满足条件属性  $B_{C1}$  且属于类  $B_{D1}$  的实例数目(考虑重复实例)。其计算方法为: 对于包含类属性的概念节点, 首先遍历其父节点, 判断父节点集合中是否存在由分类规则条件属性组成内涵的概念节点, 如果存在该类型父节点, 则求取该节点的概念外延基数, 所得规则的置信度为:  $Conf = \frac{O_D}{O_{parent}} \times 100\%$ , 其中  $O_D$  为待求节点的外延基数,  $O_{parent}$  为其对应父节点的外延基数; 如果不存在该类型父节点, 则得到的规则置信度为 100%。

## 2.4 内涵缩减

对于上述形式背景中确定的分类规则  $B_{C1} \rightarrow B_{D1}$ , 为了保证分类规则的属性无冗余性, 希望条件属性集  $C_1$  越小越好, 因此需要对概念节点在保持外延对象集不变的前提下进行内涵缩减, 具体定义为:

定义 5: 对于一个给定的概念  $C = (O_1, D_1)$ , 如果特征集合  $D_2$  满足下述两个条件:

- (1)  $g(D_2) = g(D_1) = O_1$ ;
- (2) 对于任意的  $D_3 \subset D_2$  有  $g(D_3) \supset g(D_2) = O_1$ ;

则它被称为是  $C$  的一个内涵缩减。

为了实现概念的内涵缩减, 在本研究中采用了族集最小覆盖理论, 其理论具体表示为:

定义 6: 对于一个给定的族集  $FM = \{M_1, M_2, \dots, M_n\}$ , 集合  $M$  被称为  $FM$  的最小覆盖, 如果它满足:

- (1)  $\forall M_i \in FM (M \cap M_i \neq \emptyset)$ ;

- (2)  $\forall M' \subset M (\exists M_i \in FM (M' \cap M_i \neq \emptyset))$ ;

定义 7: 对于一个概念节点  $C = (O_1, D_1)$  和一个子集  $D_2 \subseteq D_1$ ,  $D_2$  是  $C$  的一个内涵缩减当且仅当  $D_2$  是族集  $\{D_1 - D_3 \mid (O_3, D_3) \text{ 是 } C \text{ 的一个父节点}\}$  的一个最小覆盖。

通过上述定义可知, 计算概念内涵缩减的问题被转化为计算相应族集的最小覆盖集的问题, 具体计算族集的最小覆盖集的实现算法详见文献谢志鹏 (2001)。

## 3 实验与分析

选取湖北省房县作为研究区, 收集的原始数据资料中遥感影像为 Landsat ETM, 成像时间为 2000-06-15, 其他基础的地理数据包括 1:10 万数字地形图、行政区划图、2000 年土地利用类型图。

### 3.1 影像特征提取

光谱特征为 Landsat ETM 影像各波段对应的 DN 值(波段 TM6 除外, 由于该波段分辨率相对其他波段较低)。波段相关性特征选用波段间归一化差值指数, 具体计算公式为:  $NDI_{ij} = \frac{TM_i - TM_j}{TM_i + TM_j} (i, j = 1, 2, 3, 4, 5, 7 \text{ 且 } i \neq j)$ 。纹理特征以 Landsat ETM 全色波段作为提取对象, 采用灰度共生矩阵法进行纹理特征提取, 其移动窗口为  $3 \times 3$ , 移动步长为 1, 移动方向取  $45^\circ$ , 选择的纹理特征量为: 均值, 方差, 逆方差, 对比度, 非相似度, 熵, 角二阶矩和相关度。

对于基础地理数据, 首先对 1:10 万地形图数字矢量化后形成 DEM, 然后在 Arcgis9.2 平台中生成坡度图, 并以遥感影像作为参考图像, 经过重采样、空间配准处理后, 将高程、坡度作为 2 个新的“波段”特征数据。

### 3.2 数据处理与分析

为了减少人工挑选样本数据带来的误差, 试验中以 2000 年该地区土地利用类型图作为参考, 以 ENVI 4.2 开发平台实现土地利用类型图中的小斑为单位的遥感影像裁减, 为了避免因为土地利用类型图解译精度造成的数据误差, 对裁减各波段特征数据进行直方图统计, 将各波段小斑范围内出现频率次数最多的值作为该波段的特征值, 形成多源空间数据库, 其结构如图 2。

根据对研究区的实地调查统计, 研究区的地物类型主要分为: 草地、林地、旱地、水田和水体, 为

了保证分类规则提取的准确性,从数据库中选取5133条样本数据,每种类型的样本数量大致相同。为了便于比较分析,将各特征值按照公式 $Y = \frac{y - y_{min}}{y_{max} - y_{min}} \times 255$ 归一化至0—255范围内,其中波段相关性特征值的取值范围为(-1, 1),高程特征值取值范围为(200m, 2200m),坡度特征值取值范围为(0°, 69°),并依据特征值制作各地物的光谱特征曲线、波段相关性特征曲线、纹理特征曲线和波段高程曲线(图3)。

从光谱特征曲线(图3(a))可知:各地物在TM3波段上区分明显,在TM1、TM5上差异不大,水体与其他地物除TM5外其他波段上有明显区分,旱地、草地、水田、林地之间在各波段上亮度值比较近似,这是由于遥感影像季相为夏季,旱地与水田的作物生长茂盛,与草地、林地之间相似度较高,因此亮度值上差异不够显著。

图3(b)中,横坐标1—15分别代表NDI<sub>12</sub>, NDI<sub>13</sub>, NDI<sub>14</sub>, NDI<sub>15</sub>, NDI<sub>17</sub>, NDI<sub>23</sub>, NDI<sub>24</sub>, NDI<sub>25</sub>, NDI<sub>27</sub>, NDI<sub>34</sub>, NDI<sub>35</sub>, NDI<sub>37</sub>, NDI<sub>45</sub>, NDI<sub>47</sub>, NDI<sub>57</sub>。从波段间

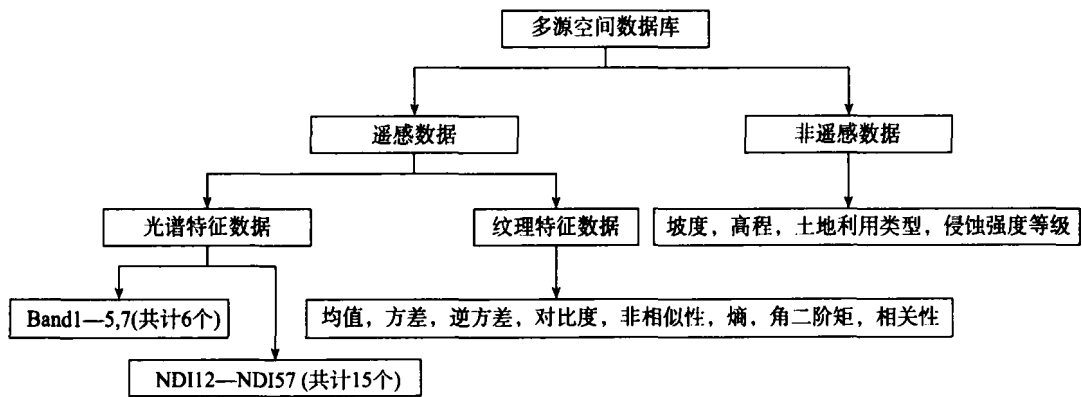


图2 多源空间数据库结构

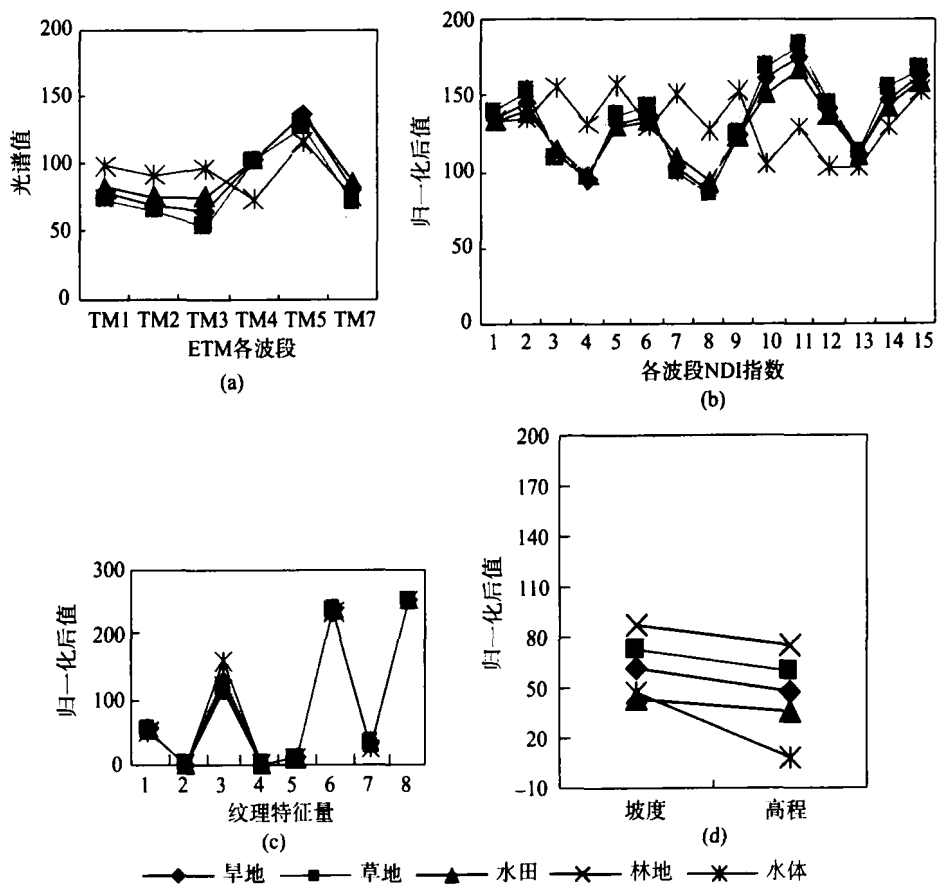


图3 研究区各地物特征曲线图

(a) 光谱特征曲线; (b) 波段间归一化差值指数曲线; (c) 纹理特征曲线; (d) 坡度高程特征曲线

归一化差值指数曲线(图 3(b))分析中可知：波段 TM3 与 TM4、TM7 与 TM4、TM3 与 TM5、TM1 与 TM3、TM1 与 TM7 之间的归一化差值指数在区分地物能力较强，而波段 TM1、TM2 之间归一化差值指数几乎不能区分地物，因此该指数在进行规则挖掘时不予选用。

图 3(c)中，横坐标 1—15 分别代表均值，方差，逆方差，对比度，非相似度，熵，角二阶矩，相关度。从纹理特征曲线分析中可知：各地物除了在逆方差特征值上有较大差异外，其他特征量几乎一致。因此在研究中，除逆方差特征量外，其他纹理特征量不予选用。

从坡度高程特征曲线分析中可知：水田、水体主要分布在高程较低、坡度较小的平原地带，旱地主要分布在坡度较高、高程相对较高的丘陵地带，海拔较高、坡度较陡的地带因受人为因素破坏较小，一般为林地或者草地。

3.3 分类规则提取及分析

由于概念格结构中只能处理离散数据，研究中依据数据的直方图分布规律，将各特征值划分为多个离散区间。对于遥感影像光谱值、波段间归一化差值指数以及纹理特征量，以 10 为间隔单位将属性空间划分 25 个区间；对于坡度属性，按照国家标准划分为 7 个区间，具体如表 2；对于高程属性，由于各地物在地形垂直带上分布一般相差 120m 左右，因此将高程特征以 120m 为间隔单位划分为 12 个区间。每个属性区间具体表示为：Name\_XX，Name 为属性名称，XX 为该属性对应的区间段。

对于离散化的样本数据，利用概念格算法生成对应的概念格结构，在该概念格结构中共产生 9033 个概念，其中包含类别属性的概念为 4835 个，经过概念内涵缩减算法处理后，得到对应的分类规则，将支持度较高的一维条件属性分类规则列出(表 3)。

表 2 坡度属性区间表

0°—3°	3°—5°	5°—8°	8°—15°	15°—25°	25°—35°	>35°
Slope_1	Slope_2	Slope_3	Slope_4	Slope_5	Slope_6	Slope_7

表 3 一维条件属性分类规则

序号	分类规则	支持度/%	置信度/%
1	band2_6—>旱地	9.71	41.88388
2	ND47_16—>草地	7.36	44.59759
3	band6_6 —>草地	6.49	40.80674
4	ND23_15—>草地	6.55	50.1248
5	ND17_11—>水田	7.51	50.55545
6	ND47_11—>水田	8.16	46.58692
7	ND47_14—>水田	4.33	55.43247
8	band4_8 —>水田	5.25	44.76058
9	ND35_14 —>水田	4.59	55.01782
10	ND13_12 —>水田	3.42	50.21281
11	ND34_11 —>水田	3.25	47.71157
12	ND35_15—>水田	5.16	49.90062
13	ND47_18—>林地	8.97	42.78857
14	ND35_17—>林地	7.35	44.75626
15	ND27_13—>林地	5.15	55.83393
16	dem_4—>林地	5.71	53.49792
17	slope_3—>水体	6.71	43.30056
18	ND14_21—>水体	7.45	44.13420
19	ND27_20—>水体	4.70	42.42167
20	band4_1 —>水体	5.44	48.22366
21	dem_0 —>水体	16.30	48.21919

从挖掘出的分类规则中可知：较高维的分类规则一般是由低维的规则组成，但是低维的分类规则具有较高的支持度，因此低维的分类规则更反映地物之间的本质区分特征。在对挖掘出的一维属性分类规则的分析中可知：对于规则 16，可以理解为当高程为 827—984m 时，地物为林地的支持度为 5.71%，可信度为 53.50%；对于规则 21，可以解释为当高程为 200—307m 时，该地物为水体的支持度为 16.30%，可信度为 48.2%，这些规则与地物的高程坡度分析结果相吻合，水体分布在地势较为平缓的平原、丘陵地带，林地分布在在地势较高的山区，因此挖掘的结果可以被理解和接受。对于规则 2, 6, 7, 13, 可知在林地波段 4 与 7 间的归一化差值指数值 ND47\_18>草地的 ND47\_16>水田的 ND47\_11—ND47\_14，这与地物在波段间相关性分析时得到的结论一致，即在波段 4 与 7 间的归一化差值指数中林地、草地、水田具有较大的区分性，且指数值林地>草地>水田，因此挖掘的结果可信度较高。

在挖掘出的分类规则长度比较上：满足支持度的分类规则长度最高由 8 个属性组成，一般为 3—6 个属性，相对于基于决策树算法 C4.5 得到的分类规则平均长度为 8—10(陈小瑜, 2007)，因此基于形式

概念分析方法所得的分类规则相对简短,以此为基础构建的分类器结构也将相对简单。

3.4 分类器构建及分析

为了实现基于挖掘出的分类规则构建分类器,采用了一种启发式分类器构造思想:首先将分类规则集合中矛盾规则予以删除;其次对规则集合按照偏序关系进行排序,排序原则为:(1)置信度优先;(2)置信度相等情况下,支持度优先;(3)置信度与支持度相等情况下,规则生成先后顺序;(4)从分类规则中挑选出对训练样本集训练错误数最少的子集构成分类器,该算法的具体实现过程详见文献齐红(2005)。为了考察不同支持度和置信度阈值对分类器精度的影响及评价分类器的分类性能,从多源空间

数据库中挑选出训练样本数据,选取不同的支持度和置信度阈值分别进行试验,试验结果表明:当满足支持度阈值为2%和置信度阈值为40%时,启发式分类器的分类精度最高,达到80.26%,且高于同一测试集下C4.5算法的分类精度78.64%。

为了验证与比较基于形式概念分析理论挖掘出的分类规则构建的启发式分类器的精度,本文选取2000-05-05轨道号为126/38的遥感影像为实验对象,以已解译的土地利用类型图与地形数据等作为辅助数据,在实验区的影像上选择检测样本1166个,其中旱地275个,草地107个,水田236个,水体228个,林地320个,基于混淆矩阵方法进行精度评价,并与监督分类法和基于决策树分类方法进行对比,结果如表4—表6。

表 4 试验区监督分类精度评价结果

类型	旱地	草地	水田	水体	林地	总和	用户精度/%
旱地	204	32	25	0	14	275	74.18
草地	19	62	17	0	9	107	57.94
水田	33	14	169	8	12	236	71.61
水体	0	0	25	203	0	228	89.04
林地	23	38	25	9	225	320	70.31
总和	279	146	261	220	260	1166	—
生产者精度/%	73.12	42.47	64.75	92.27	86.54	—	—
总精度: 74.01% Kappa 系数: 0.7389							

表 5 试验区决策树 C4.5 分类精度评价结果

类型	旱地	草地	水田	水体	林地	总和	用户精度/%
旱地	215	21	28	1	10	275	78.18
草地	21	71	8	0	7	107	66.35
水田	35	9	173	4	15	236	73.31
水体	2	5	8	205	8	228	89.91
林地	20	42	19	5	234	320	73.13
总和	293	148	236	215	274	1166	—
生产者精度/%	73.38	47.97	73.31	95.35	85.40	—	—
总精度: 77.02% Kappa 系数: 0.7691							

表 6 试验区基于构建的启发式分类器分类精度评价结果

类型	旱地	草地	水田	水体	林地	总和	用户精度/%
旱地	221	25	14	2	13	275	80.36
草地	23	75	3	1	5	107	70.09
水田	27	12	181	2	14	236	76.69
水体	4	6	15	199	4	228	87.28
林地	18	37	17	6	242	320	75.63
总和	293	155	230	210	278	1166	—
生产者精度/%	75.43	48.39	78.70	94.76	87.05	—	—
总精度: 78.73% Kappa 系数: 0.7862							

从表4—表6的评价结果可以看出: 基于形式概念分析挖掘出的分类规则构建的分类器与决策树 C4.5 算法和监督分类法相比, 在分类总精度和 Kappa 系数比较中相对较高, 在用户精度与生产者精度的比较中, 也表现出一定的优势, 特别是旱地、草地、水田、林地的区分精度相对监督分类法有较大提高, 与决策树 C4.5 算法相比也有一定的提高, 但是对水体的分类精度略有下降。由于缺乏实地观测数据, 实验中选取的参考数据为解译该区土地利用类型图有一定的精度误差, 造成 3 种分类方法的总体分类精度相对较低, 但是从 3 种方法之间的比较结果看, 基于形式概念分析挖掘出的分类规则构建的分类器还是体现了一定的优势, 因此证明了该方法对遥感影像分类提供了一种新的思路, 对解决遥感影像的“同谱异物、同物异谱”现象具有一定的借鉴意义。

## 4 结论与展望

本文以形式概念分析理论为基础, 实现了该理论在遥感图像分类规则挖掘中的应用, 并以挖掘出的分类规则为基础实现了启发式分类器的构造, 实验结果表明该分类方法与监督分类法、决策树 C4.5 算法相比表现出一定的优势, 证明了该方法为遥感影像的分类规则提取提供了一种新的研究思路与方法。但是, 鉴于形式概念分析理论在挖掘分类规则时对样本数据噪声敏感的特点, 在处理遥感影像中地物类型混杂的区域时, 需考虑混合像元分解, 因此如何利用形式概念分析理论实现混合像元分解有待进一步研究。

## REFERENCES

- Chen X Y. 2007. The Classification of Remote Sensing Image Based on Spatial Data Mining and Knowledge Discovery. Fujian: FuJian Normal University
- Chen Y F, Wang Z Q, Zhang Y G and Zhang Y Z. 1996. Expert system applied for TM image classification. *Forest Research*, 4(9): 344—347
- He L M. 2006. Support Vector Machines Ensemble and Its Application in Remote Sensing Classification. Zhejiang: Zhejiang

University

- Jian S Q, Hu X G and Jiang M H. 2001. Incremental algorithms of extended concept lattice. *Computer Engineering and Applications*, 15 (37): 132—134
- Qi H. 2005. Research on Knowledge Discovery based on Lattice Concept Analysis. Jiling: Jiling University
- Sun J B, Shu N and Guan Z Q. 1997. The Remote Sensing Theory, Method and Application. Beijing: Surveying and Mapping Press
- Sun X B, Fan W, Yan P, Huang Y and Ma Y H. 2007. A brief review on the classification methods of land cover based on remote sensing image. *Chinese Agricultural Science Bulletin*, 9(23): 607—610
- Wille R. 1989. Knowledge Acquisition by Methods of Formal Concept Analysis, Learning Symbolic and Numeric Knowledge. New York: Nova Science Publishers, Inc
- Wille R. 1982. Restructuring lattice theory: an approach based on hierarchies of concepts. *Ordered Sets*. Dordrecht: D Reidel Publing Company
- Xie Z P. 2001. Research on Knowledge Discovery Based on Concept Lattice Model. Hefei: Hefei University of Technology
- Xiong Z, Tong Q X and Zhen L F. 2000. High-rank artificial neural network algorithm for classification of hyper spectral image data. *Journal of Image and Graphics*, 3(5): 196—201
- Zhai Y H. 2006. Knowledge Representation and Knowledge Acquisition on Formal Concept Analysis. Shanxi: Shanxi University

## 附中文参考文献

- 陈小瑜. 2007. 基于空间数据挖掘与知识发现的遥感影像分类研究. 福建师范大学
- 陈永富, 王振琴, 张玉贵, 张彦忠. 1996. 专家系统在 TM 遥感图像分类中的应用研究. *林业科学研究*, 4(9): 344—347
- 何灵敏. 2006. 支持向量机集成及在遥感分类中的应用. 浙江大学
- 简宋全, 胡学钢, 蒋美华. 2001. 扩展概念格的渐进式构造. *计算机工程与应用*, 15(37): 132—134
- 齐红. 2005. 基于形式概念分析的知识发现方法研究. 吉林大学
- 孙家柄, 舒宁, 关泽群. 1997. 遥感原理、方法和应用. 北京: 测绘出版社
- 孙秀邦, 范伟, 严平, 黄勇, 马友华. 2007. 遥感影像土地覆被分类研究进展. *中国农学通报*, 9(23): 607—610
- 谢志鹏. 2001. 基于概念格模型的知识发现研究. 合肥工业大学
- 熊桢, 童庆禧, 郑兰芬. 2000. 用于高光谱遥感图像分类的一种高阶神经网络算法. *中国图象图形学报*, 3(5): 196—201
- 翟岩慧. 2006. 形式概念分析中的知识表示与知识获取. 山西大学